

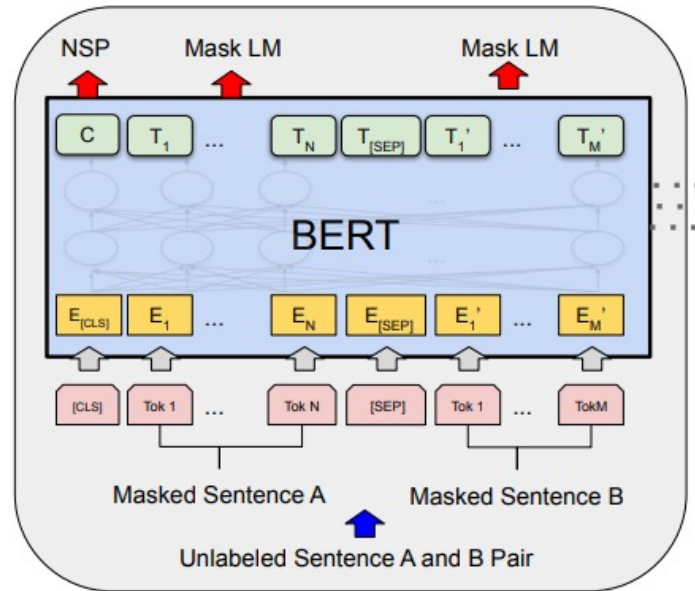
# Generative Video Transformer: Can Objects be the Words?

Yi-Fu Wu<sup>1</sup>, Jaesik Yoon<sup>1,2</sup>, Sungjin Ahn<sup>1</sup>

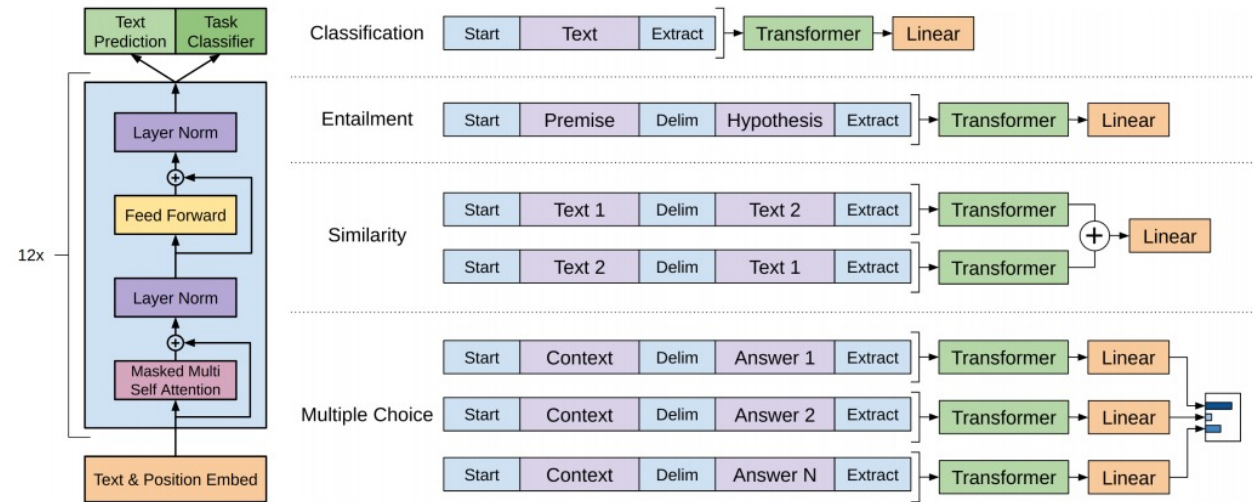
<sup>1</sup>Rutgers University, <sup>2</sup>SAP



# Motivation



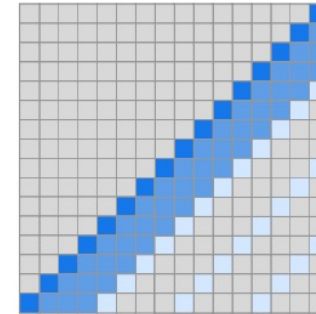
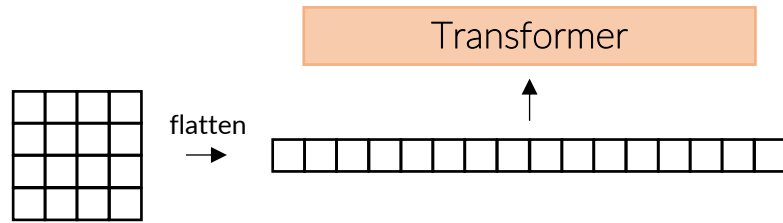
BERT (Devlin et al., 2019)



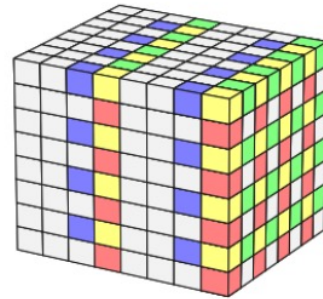
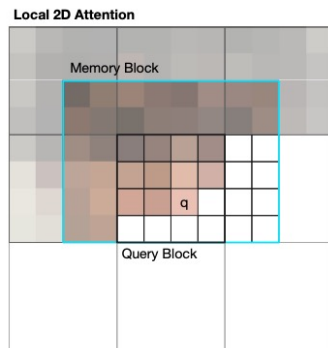
GPT (Radford et al., 2018)

- Can a transformer-based architecture be effective for generative pre-training of visual scenes for video generation and understanding?

# Previous Work

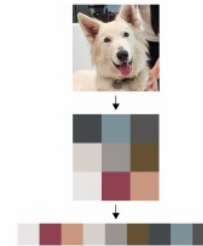


*Child et al., 2019*

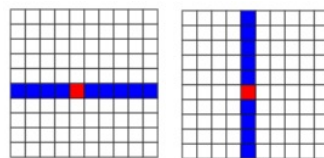


*Parmar et al., 2018*

*Parmar et al., 2018*



*Chen et al., 2020*

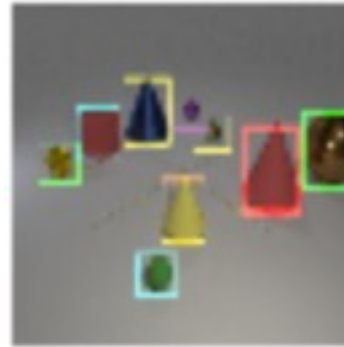
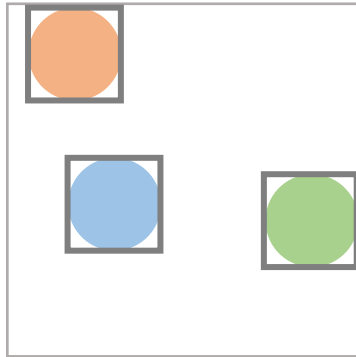


*Ho et al., 2019*



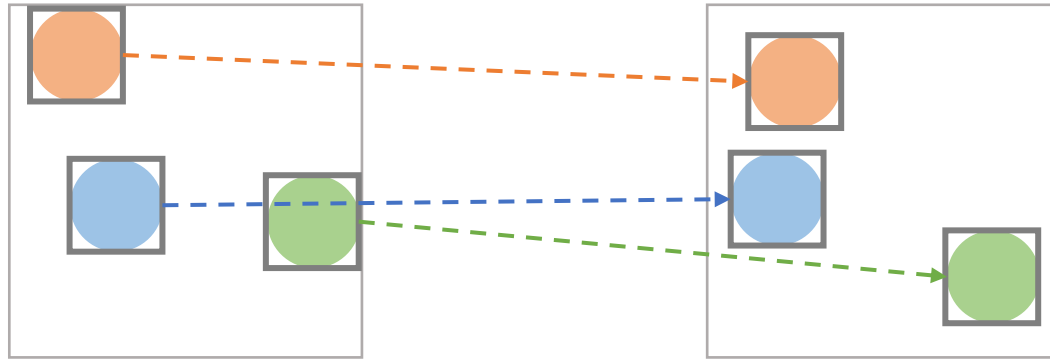
*Dosovitskiy et al., 2021*

# Our Approach



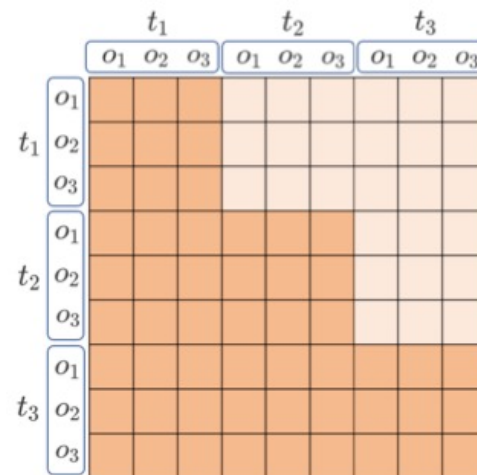
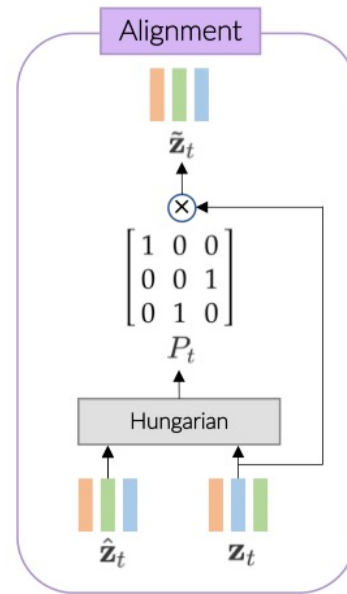
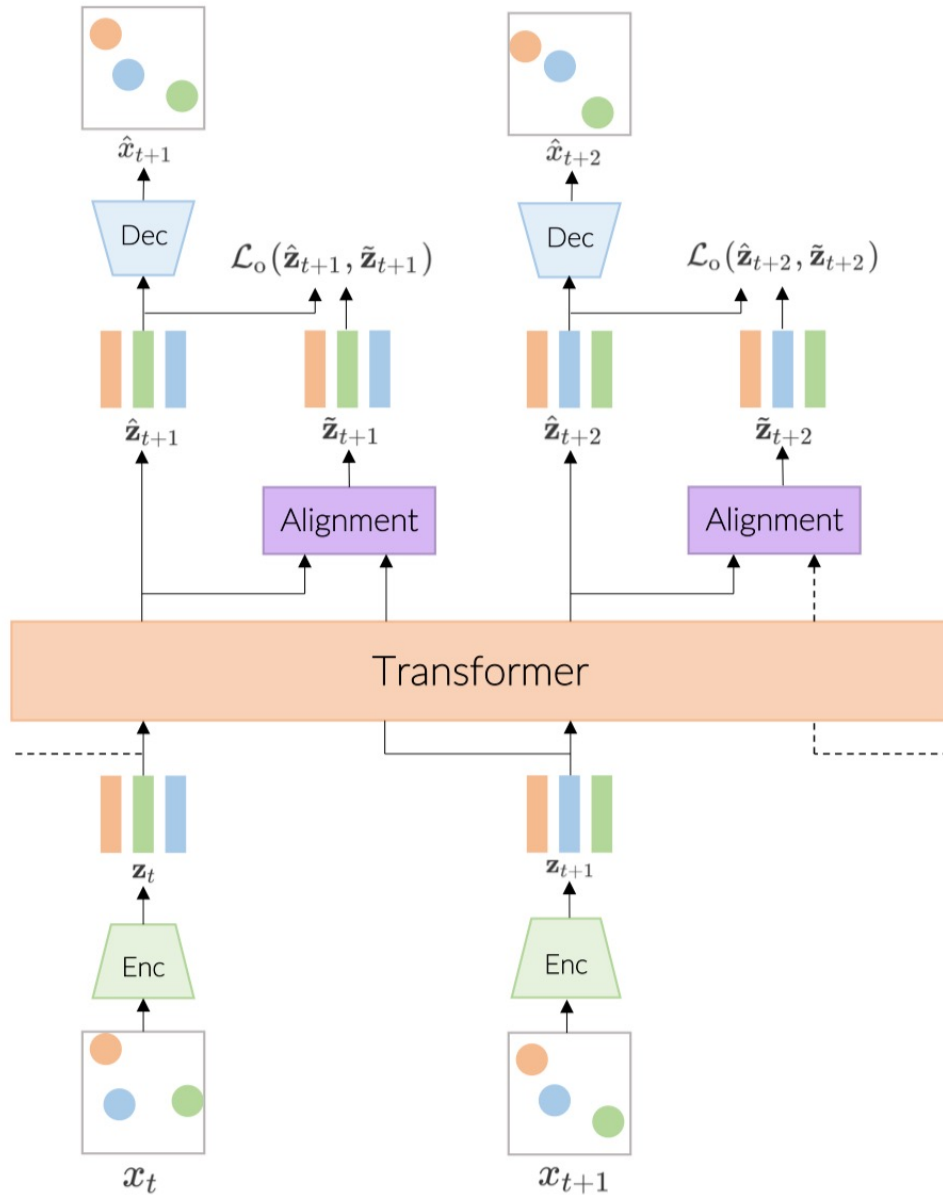
- Tokenize an image into its constituent object representations and use these representations as input to the transformer

# Our Approach



- Predict the entire image at once by generating all the objects in an image simultaneously given their previous states
- Align objects between frames based on object location, leveraging SPACE (*Lin et al., 2020*) an unsupervised object representation model that outputs explicit bounding box information

# Object-Centric Video Transformer

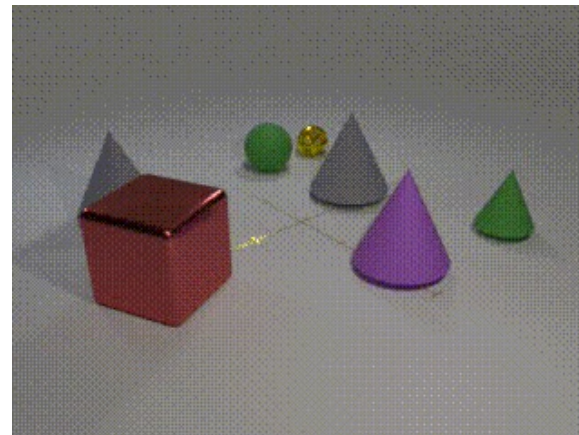
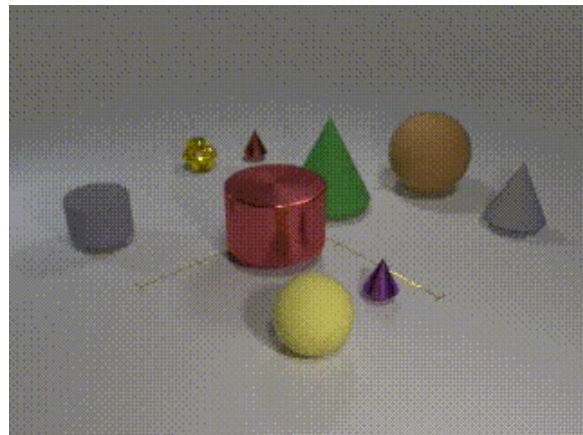


# Experiments - Datasets

- Bouncing balls



- CATER



# Experiments – Bouncing Ball

Table 1. Average next-step prediction color change accuracy

	MOD1	MOD2	MOD3	MOD1234
GSWM	71.69	17.51	14.63	11.72
LSTM+GNN	73.64	69.08	22.30	51.38
SVVT	37.53	18.23	11.96	29.47
CONVVT	88.31	82.83	46.49	67.29
CONVVT-AR	8.70	4.20	3.25	6.10
OCVT-AR	78.70	76.99	54.49	64.97
<b>OCVT (OURS)</b>	<b>89.61</b>	<b>88.18</b>	<b>82.70</b>	<b>78.43</b>

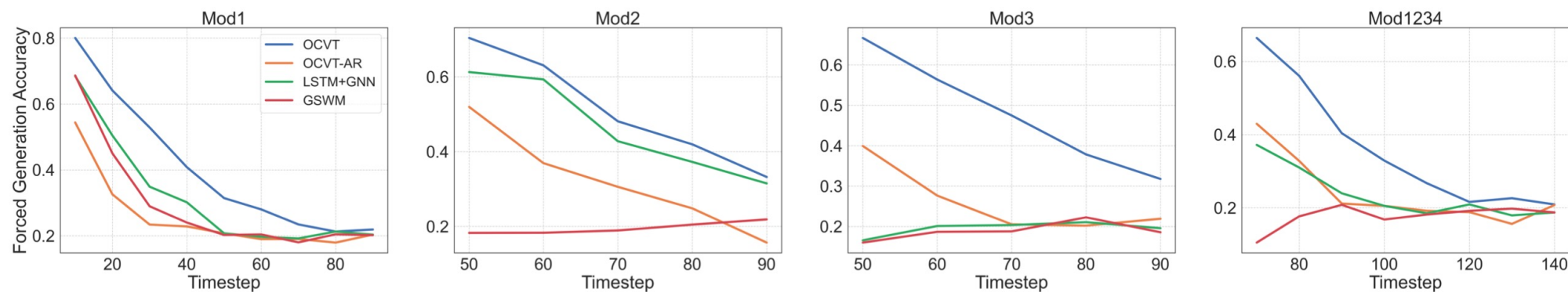
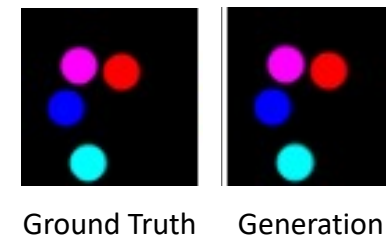


Figure 5. Forced Generation Accuracy



# Experiments - CATER

*Table 2. CATER results*

	TOP 1 $\uparrow$	TOP 5 $\uparrow$	L1 $\downarrow$
DING ET AL	70.6	93.0	0.53
DING ET AL W/ L1	74.0	94.0	0.44
HOPPER	73.2	93.8	0.85
OPNET	74.8	-	0.54
OCVT (OURS)	<b>76.0</b>	94.4	0.45
OCVT W/ L1 (OURS)	75.9	<b>95.3</b>	<b>0.39</b>

# Conclusion

- OCVT is able to generate future frames of videos with complex long-term dependencies
- Learned representations are useful for downstream tasks
- *Please refer to our paper for more details*